

(19) 日本国特許庁(JP)

(12) 公開特許公報(A)

(11) 特許出願公開番号

特開2008-257511

(P2008-257511A)

(43) 公開日 平成20年10月23日(2008.10.23)

(51) Int.Cl.	F I	テーマコード (参考)
<b>G06F 17/28 (2006.01)</b>	G06F 17/28 U	5B075
<b>G06F 17/30 (2006.01)</b>	G06F 17/30 210A	5B091

審査請求 未請求 請求項の数 10 O L (全 24 頁)

(21) 出願番号 特願2007-99403 (P2007-99403)  
 (22) 出願日 平成19年4月5日(2007.4.5)

(71) 出願人 500257300  
 ヤフー株式会社  
 東京都港区六本木六丁目10番1号  
 (74) 代理人 100106002  
 弁理士 正林 真之  
 (72) 発明者 萩原 健  
 東京都港区六本木6丁目10番1号 ヤフー株式会社内  
 (72) 発明者 増山 毅司  
 東京都港区六本木6丁目10番1号 ヤフー株式会社内  
 (72) 発明者 本野 秀樹  
 東京都港区六本木6丁目10番1号 ヤフー株式会社内  
 Fターム(参考) 5B075 NK32

最終頁に続く

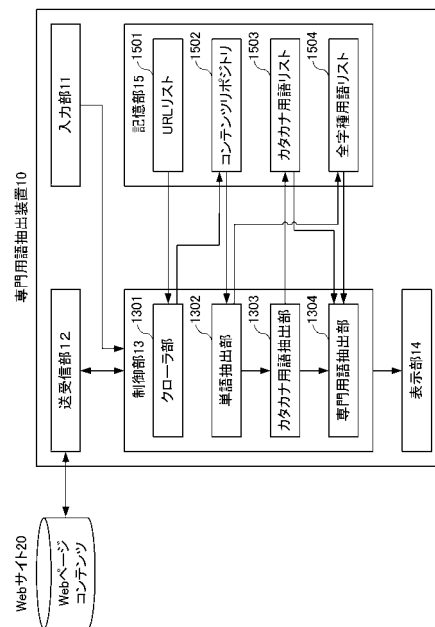
(54) 【発明の名称】 専門用語抽出装置、方法及びプログラム

(57) 【要約】

【課題】 Webドキュメントから専門用語を自動的に抽出する専門用語抽出装置を提供すること。

【解決手段】 本装置のクロール部が、専門分野ごとのURLリストを用いて、Webページのコンテンツを収集する。URLリストは、常にUp-To-Dateに更新する。次に、本装置の単語抽出部によって、収集されたWebページのコンテンツのテキストを形態素解析し品詞に分類して、カタカナ語彙と全字種の語彙を抽出する。この際、助詞や接続詞など専門用語になりにくい品詞は抽出対象から除外する。そして、本装置のカタカナ用語抽出部によって、抽出されたカタカナ語彙に対して、FLR法を用いて、重要度の計算を行い重要度の高いカタカナ用語を抽出する。さらに、専門用語抽出部によって、カタカナ用語と、先に抽出された全字種の語彙との共起ヒット情報を計算して、専門用語を抽出する。

【選択図】 図2



**【特許請求の範囲】****【請求項 1】**

Web ページから専門用語を抽出する専門用語抽出装置であって、  
専門分野ごとに定められた URL リストに含まれた URL にアクセスし、前記 Web ページのコンテンツを収集するクローラ部と、  
前記コンテンツを形態素解析し、カタカナ語彙と全字種語彙を抽出する単語抽出部と、  
前記抽出されたカタカナ語彙に対して、FLR 法を用いて前記カタカナ語彙からカタカナ用語を抽出するカタカナ用語抽出部と、  
前記カタカナ用語抽出部によって抽出されたカタカナ用語と、前記単語抽出部によって抽出された全字種語彙とを、TFIDF 値とシン普森係数値を組み合わせた共起ヒット情報を計算して、専門用語を抽出する専門用語抽出部と、  
を備えた専門用語抽出装置。

10

**【請求項 2】**

前記カタカナ用語抽出部は、前記 FLR 法に代えて、C-Value 法を用いる、請求項 1 に記載の装置。

**【請求項 3】**

前記カタカナ用語抽出部は、前記 FLR 法に代えて、MC-Value 法を用いる、請求項 1 に記載の装置。

**【請求項 4】**

前記専門用語抽出部は、前記シン普森係数値に代えて、相互情報量値を用いる、請求項 1 乃至 3 に記載の装置。

20

**【請求項 5】**

前記専門用語抽出部は、前記シン普森係数値に代えて、ダイス係数値を用いる、請求項 1 乃至 3 に記載の装置。

**【請求項 6】**

前記専門用語抽出部は、前記シン普森係数値に代えて、ジャガード係数値を用いる、請求項 1 乃至 3 に記載の装置。

**【請求項 7】**

前記専門用語抽出部は、前記シン普森係数値に代えて、コサイン類似度値を用いる、請求項 1 乃至 3 に記載の装置。

30

**【請求項 8】**

前記専門用語として、アダルト専門分野における掲載禁止用語を抽出する請求項 1 乃至 7 に記載の装置。

**【請求項 9】**

Web ページから専門用語を抽出するための方法であって、  
専門分野ごとに定められた URL リストに含まれた URL にアクセスし、前記 Web ページのコンテンツを収集するステップと、  
前記コンテンツを形態素解析し、カタカナ語彙と全字種語彙を抽出する単語抽出ステップと、  
前記抽出されたカタカナ語彙に対して、FLR 法を用いて前記カタカナ語彙からカタカナ用語を抽出するカタカナ用語抽出ステップと、  
前記カタカナ用語抽出ステップによって抽出されたカタカナ用語と、前記単語抽出ステップによって抽出された全字種語彙とを、TFIDF 値とシン普森係数値を組み合わせた共起ヒット情報を計算して、専門用語を抽出する専門用語抽出ステップと、  
を含む方法。

40

**【請求項 10】**

Web ページから専門用語を抽出するためのコンピュータ・プログラムであって、コンピュータに、  
専門分野ごとに定められた URL リストに含まれた URL にアクセスし、前記 Web ページのコンテンツを収集するステップと、

50

前記コンテンツを形態素解析し、カタカナ語彙と全字種語彙を抽出する単語抽出ステップと、

前記抽出されたカタカナ語彙に対して、FLR法を用いて前記カタカナ語彙からカタカナ用語を抽出するカタカナ用語抽出ステップと、

前記カタカナ用語抽出ステップによって抽出されたカタカナ用語と、前記単語抽出ステップによって抽出された全字種語彙とを、TFIDF値とシン普森係数値を組み合わせた共起ヒット情報を計算して、専門用語を抽出する専門用語抽出ステップと、

を実行させるコンピュータ・プログラム。

【発明の詳細な説明】

【技術分野】

10

【0001】

本発明は、ドキュメントからの専門用語抽出装置、方法及びプログラムに関する。特に、Webドキュメントからの専門用語抽出装置、方法及びプログラムに関する。

【背景技術】

【0002】

様々な分野において、専門用語のデータベースを作成したり、データベースから専門用語を検索したりするために、専門分野のドキュメントから専門用語の抽出が行われている。従来、専門用語の抽出は当該分野の専門家が人手でドキュメントを精査し、抽出していたが、その作業を自動化するための試みが複数なされている。例えば、非特許文献1には、単名詞を含む単名詞パイグラムの左右に接続する単名詞を抽出し、その頻度を基にスコアリングを行い、専門用語を抽出する方法が開示されている。又、特許文献1には、大量の専門用語が抽出される分野において、専門用語辞書を最新状態にメンテナンスするために、ある用語の関連語の同族語、類似語の同族語を抽出することで、多様な周辺語彙を網羅的に情報収集し、新語登録などのメンテナンス作業を効率化する方法が開示されている。

20

【非特許文献1】出現頻度と接続頻度に基づく専門用語抽出、湯本他、自然言語処理、10(1)27-45, 2003年1月

【特許文献1】特開2005-222263号公報

【発明の開示】

【発明が解決しようとする課題】

30

【0003】

しかしながら、特許文献1及び非特許文献1に記載の技術では共に、専門用語を抽出する対象ドキュメントが既にデータベースに保存されており、かつ専門用語と関連する分野のドキュメントである(特許文献1であれば、医学・生物分野、非特許文献1であれば、情報処理分野)ことを前提としている。そのため、対象とするドキュメント数が限定されて、高精度で専門用語を抽出することができた。しかし、対象をWebサイト全体に広げた場合、Webドキュメントは分野ごとに分類されていないという問題があり、専門用語を抽出する前に、まず対象とするWebドキュメントをWeb上から収集する必要がある。又、Webサイトは次々に更新されるという特徴があり、さらに企業や官公庁だけでなく、個人の趣味・嗜好の基に作成されるものも多く存在するため、学术论文などに比べてノイズとなる情報がドキュメント中に多く含まれている可能性が高く、上記の技術とは別の視点が必要となる。

40

【0004】

本発明は、上記課題に鑑み、Webドキュメントから専門用語を自動的に抽出する専門用語抽出装置を提供することを目的とする。

【課題を解決するための手段】

【0005】

本発明では以下のような解決手段を提供する。

【0006】

(1) Webページから専門用語を抽出する専門用語抽出装置であって、

50

専門分野ごとに定められたURLリストに含まれたURLにアクセスし、前記Webページのコンテンツを収集するクローラ部と、

前記コンテンツを形態素解析し、カタカナ語彙と全字種語彙を抽出する単語抽出部と、

前記抽出されたカタカナ語彙に対して、FLR法を用いて前記カタカナ語彙からカタカナ用語を抽出するカタカナ用語抽出部と、

前記カタカナ用語抽出部によって抽出されたカタカナ用語と、前記単語抽出部によって抽出された全字種語彙とを、TFIDF値とシンブソン係数値を組み合わせた共起ヒット情報を計算して、専門用語を抽出する専門用語抽出部と、

を備えた専門用語抽出装置。

#### 【0007】

10

(1)の構成によれば、まず、本装置に備えられたクローラ部が、専門分野ごとに分けられたURLリストを用いて、Webページのコンテンツを収集(クローラ)する。URLリストは、常にUp-Dateに更新する。次に、本装置の単語抽出部によって、収集されたWebページのコンテンツのテキストを形態素解析し品詞に分類して、カタカナ語彙と全字種の語彙を抽出する。この際、助詞や接続詞など専門用語になりにくい品詞は抽出対象から除外してよい。そして、本装置のカタカナ用語抽出部によって、カタカナ語彙から、FLR法を用いて、重要度の計算を行い重要度の高いカタカナ用語を抽出する。さらに、専門用語抽出部によって、抽出されたカタカナ用語と、先に抽出された全字種の語彙との共起ヒット情報(2つの語彙が同じドキュメントで共起する度合い)を計算することによって専門用語を抽出する。なお、FLR法とは、後述するように、接続頻度LR法(接続種類LR法)に、用語Wがコーパス(言語資料体)中に出現した頻度を加味したものである。

20

#### 【0008】

このように、まずカタカナ語彙に着目して重要度の高いカタカナ用語を求めるカタカナ用語抽出処理と、この重要度の高いカタカナ用語と全字種の語彙との共起ヒット情報による専門用語抽出処理を行うことによって、特にカタカナ語彙を含んだ専門用語(医薬分野、IT分野、ロボット工学分野、アダルト分野などの専門用語)に対して、膨大に存在するWebドキュメントから、Up-to-Dateに専門用語を自動的に抽出することが可能になる。

#### 【0009】

30

(2) 前記カタカナ用語抽出部は、前記FLR法に代えて、C-Value法を用いる、(1)に記載の装置。C-Value法は、後述するように、用語Wを部分文字列として含むより長い用語の出現頻度を、用語Wを部分文字列として含むより長い用語の種類数で割った値を用語Wの出現頻度から補正した値を重要度とする方法である。

#### 【0010】

(2)の構成によれば、カタカナ用語抽出部において、FLR法に代えて公知のC-Value法を用いることができる。

#### 【0011】

(3) 前記カタカナ用語抽出部は、前記FLR法に代えて、MC-Value法を用いる、(1)に記載の装置。

40

#### 【0012】

(3)の構成によれば、カタカナ用語抽出部において、FLR法に代えてC-Value法を改良したMC-Value法(Modified C-Value法)を用いることができる。

#### 【0013】

(4) 前記専門用語抽出部は、前記シンブソン係数値に代えて、相互情報量値を用いる、(1)乃至(3)に記載の装置。

#### 【0014】

(5) 前記専門用語抽出部は、前記シンブソン係数値に代えて、ダイス係数値を用いる、(1)乃至(3)に記載の装置。

50

## 【 0 0 1 5 】

( 6 ) 前記専門用語抽出部は、前記シンプソン係数値に代えて、ジャガード係数値を用いる、( 1 )乃至( 3 )に記載の装置。

## 【 0 0 1 6 】

( 7 ) 前記専門用語抽出部は、前記シンプソン係数値に代えて、コサイン類似度値を用いる、( 1 )乃至( 3 )に記載の装置。

## 【 0 0 1 7 】

( 4 )から( 7 )の構成によれば、専門用語抽出部において、TFIDF法とシンプソン係数を組み合わせた方法以外にも共起ヒット情報を求める各種の手段(相互情報量値、ダイス係数値、ジャガード係数値、コサイン類似度値)を活用することができる。

10

## 【 0 0 1 8 】

( 8 ) 前記専門用語として、アダルト専門分野における掲載禁止用語を抽出する( 1 )乃至( 7 )に記載の装置。

## 【 0 0 1 9 】

( 8 )の構成によれば、専門分野としてペアレンタルコントロールに着目し、有害サイト、特にアダルトサイトで使用されるような「掲載禁止用語」(以下、NG語彙とも呼ぶ)を抽出する。アダルトサイトは規制しても次々と新しいサイトが出現し、又NG語彙にはカタカナが多く使用されるので、このようなNG語彙を含んだサイトのフィルタリングに本発明の手法が有効である。

## 【 0 0 2 0 】

( 9 ) Webページから専門用語を抽出するための方法であって、  
専門分野ごとに定められたURLリストに含まれたURLにアクセスし、前記Webページのコンテンツを収集するステップと、

20

前記コンテンツを形態素解析し、カタカナ語彙と全字種語彙を抽出する単語抽出ステップと、

前記抽出されたカタカナ語彙に対して、FLR法を用いて前記カタカナ語彙からカタカナ用語を抽出するカタカナ用語抽出ステップと、

前記カタカナ用語抽出ステップによって抽出されたカタカナ用語と、前記単語抽出ステップによって抽出された全字種語彙とを、TFIDF値とシンプソン係数値を組み合わせた共起ヒット情報を計算して、専門用語を抽出する専門用語抽出ステップと、

30

を含む方法。

## 【 0 0 2 1 】

( 9 )の構成によれば、( 1 )と同様の作用効果を持つ発明を方法として提供できる。

## 【 0 0 2 2 】

( 1 0 ) Webページから専門用語を抽出するためのコンピュータ・プログラムであって、

コンピュータに、

専門分野ごとに定められたURLリストに含まれたURLにアクセスし、前記Webページのコンテンツを収集するステップと、

前記コンテンツを形態素解析し、カタカナ語彙と全字種語彙を抽出する単語抽出ステップと、

40

前記抽出されたカタカナ語彙に対して、FLR法を用いて前記カタカナ語彙からカタカナ用語を抽出するカタカナ用語抽出ステップと、

前記カタカナ用語抽出ステップによって抽出されたカタカナ用語と、前記単語抽出ステップによって抽出された全字種語彙とを、TFIDF値とシンプソン係数値を組み合わせた共起ヒット情報を計算して、専門用語を抽出する専門用語抽出ステップと、

を実行させるコンピュータ・プログラム。

## 【 0 0 2 3 】

( 1 0 )の構成によれば、( 1 )と同様の作用効果を持つ発明をコンピュータ・プログラムとして提供できる。

50

## 【発明の効果】

## 【0024】

本発明によれば、カタカナ語彙が専門用語として多く使用される専門分野において、次々と更新されるWebサイト上の膨大なドキュメント群から、人手で精査することなく自動で専門用語抽出を行うことができる。

## 【発明を実施するための最良の形態】

## 【0025】

以下、本発明の実施形態について図を参照しながら説明する。

## 【0026】

## [システムの全体構成]

図1は、本発明の一実施形態に係るシステム1の全体構成を示す図である。

## 【0027】

本実施形態におけるシステム1は、テキストや画像などを含んだドキュメントデータ（例えば、インターネットやイントラネット上のWebページ）の解析を行い、ドキュメントデータに含まれる用語群を抽出して、該用語群から専門用語の抽出を行うシステムである。本システム1では、専門用語抽出装置10が、通信ネットワーク30を介して、様々なWebサイト20と接続される。専門用語抽出装置10は、専用装置であっても、他の目的のサーバ上に実現してもよい。なお、専門用語抽出装置10のハードウェアの数に制限はなく、必要に応じて、1又は複数のハードウェアで構成してもよい。

## 【0028】

Webサイト20は、Webページを蓄積しており、通信ネットワーク30、例えば、インターネットなどのネットワークを通じて、これらの情報をユーザの端末に送信する機能を有している。なお、個人や会社のホームページなどのWebページ群、又はWebページ群が置いてあるインターネット、又はイントラネット上の場所を、Webサイトという。

## 【0029】

通信ネットワーク30は、例えば、インターネットであり、通信回線は有線により実現するものだけではなく、アクセスポイントを介して無線LANにより実現するものなど、本発明の技術的思想に合致するものであれば様々な通信技術により実現される。

## 【0030】

専門用語抽出装置10は、専門分野ごとのURLリスト1501(a、b、c、d、・・)にあるURLのWebサイト20を参照し、該Webサイト20のWebページデータ(コンテンツ)を、通信ネットワーク30を介して収集する。そして、収集したWebページデータをコンテンツリポジトリ1502に記憶する。さらに、収集したWebページに含まれるテキストデータを形態素解析して、語彙を抽出し、専門用語を抽出する機能を備える。

## 【0031】

ここで、URLリスト1501は、管理者が、特定の分野のWebサイト20のURLをリストにすることによって与えられるものとする。例えば、特定の分野とは、情報処理分野のWebサイト20(URLリスト1501a)、医療・生物分野のWebサイト20(URLリスト1501b)、アダルト専門分野のWebサイト20(URLリスト1501c)、又はロボット工学関連分野のWebサイト20(URLリスト1501d)などである。こうすることで、特定の分野における専門用語を抽出することができる。ここでは、URLリスト1501が複数ある例を示しているが、1つのURLリスト1501に、URLと特定の分野を関連付けて記憶することで実現してもよい。

## 【0032】

なお、アダルト専門分野のWebサイト20(URLリスト1501c)から専門用語を抽出するということは、公序良俗に反するような用語を抽出することである。そして、抽出した用語を掲載禁止用語(NGワード)とし、このNGワードを含むWebサイトの検索に用いたり、有害サイトの特定に用いることができる。

10

20

30

40

50

## 【 0 0 3 3 】

## [ 専門用語抽出装置の機能ブロック ]

図 2 は、本発明の一実施形態に係る専門用語抽出装置 1 0 の機能ブロック図である。

## 【 0 0 3 4 】

専門用語抽出装置 1 0 は、主として入力部 1 1、送受信部 1 2、制御部 1 3、表示部 1 4、及び記憶部 1 5 により構成される。入力部 1 1 は、キーボード及びマウスなどの入力装置を含み、専門用語抽出装置 1 0 に対する管理者などからの入力を受け付ける機能を有している。又、送受信部 1 2 は、任意の通信インターフェイスを含み、装置からリクエストを Web サイト 2 0 に送信する機能、及び Web サイト 2 0 の Web ページデータを受信する機能を有している。さらに、制御部 1 3 は、CPU ( C e n t r a l P r o c e s s i n g U n i t ) を含み、専門用語抽出装置 1 0 を制御する機能を有している。そして、表示部 1 4 は、ブラウン管表示装置 ( C R T ) や液晶ディスプレイ ( L C D ) などの表示装置を含み、データを表示する機能を有している。又さらに、記憶部 1 5 は、ハードディスクなどの内部又は外部の記憶装置を含み、データを記憶する機能を有している。

10

## 【 0 0 3 5 】

専門用語抽出装置 1 0 の制御部 1 3 は、クローラ部 1 3 0 1、単語抽出部 1 3 0 2、カタカナ用語抽出部 1 3 0 3、及び専門用語抽出部 1 3 0 4 を有している。クローラ部 1 3 0 1 は、通信ネットワーク 3 0 を介して、Web ページなどのドキュメントデータを収集する。なお、クローラとは一般的に検索ロボットともいわれ、通信ネットワーク 3 0 を通じて、Web サイト 2 0 から Web ページデータを収集するプログラムである。そして、クローラが、Web サイトを探し出す手段や、対象とする Web ページデータの種別は様々であり、クローラの管理者の設定により、収集される Web ページデータの種別や分野も異なる。

20

## 【 0 0 3 6 】

又、単語抽出部 1 3 0 2 は、ドキュメント中のテキストを形態素解析して、単語を抽出し、カタカナ語彙と、全字種の語彙とに分けて、全字種の語彙を全字種用語リスト 1 5 0 4 に記憶する。そして、カタカナ用語抽出部 1 3 0 3 は、カタカナ語彙の用語ごとに重要度 ( 後述 ) を計算し、管理者の設定する閾値以上の用語を抽出し、カタカナ用語リスト 1 5 0 3 に記憶する。さらに、専門用語抽出部 1 3 0 4 は、カタカナ用語リスト 1 5 0 3 と、全字種用語リスト 1 5 0 4 とにおいて共起の強い用語を専門用語として抽出する。

30

## 【 0 0 3 7 】

専門用語抽出装置 1 0 の記憶部 1 5 は、URL リスト 1 5 0 1、コンテンツリポジトリ 1 5 0 2、カタカナ用語リスト 1 5 0 3、及び全字種用語リスト 1 5 0 4 を含んで構成される。URL リスト 1 5 0 1 は、クローラ部 1 3 0 1 による Web ページデータ収集先の Web サイト 2 0 の URL を記憶する。又、コンテンツリポジトリ 1 5 0 2 は、クローラ部 1 3 0 1 により収集された Web ページデータを記憶する。そして、カタカナ用語リスト 1 5 0 3 は、カタカナ用語を記憶する。さらに、全字種用語リスト 1 5 0 4 は、全字種の語彙を記憶する。

## 【 0 0 3 8 】

## [ 専門用語抽出処理 ]

図 3 は、本発明の一実施形態に係る専門用語抽出処理のフローチャートである。

40

## 【 0 0 3 9 】

まず、ステップ S 1 0 1 では、専門用語抽出装置 1 0 の制御部 1 3 が、送受信部 1 2 を介して、クローラ部 1 3 0 1 により、Web ページなどのドキュメントデータを収集する。なお、記憶部 1 5 の URL リスト 1 5 0 1 に含まれた URL に対する、Web サイト 2 0 の Web ページデータを収集してもよい。

## 【 0 0 4 0 】

次に、ステップ S 1 0 2 では、クローラ部 1 3 0 1 が、収集したドキュメントデータを、コンテンツリポジトリ 1 5 0 2 に記憶する。

## 【 0 0 4 1 】

50

次に、ステップ S 1 0 3 では、単語抽出部 1 3 0 2 が、コンテンツリポジトリ 1 5 0 2 から、ドキュメントデータを読み込む。

【 0 0 4 2 】

次に、ステップ S 1 0 4 では、単語抽出部 1 3 0 2 が、ドキュメントデータのテキストを形態素解析する。ここで、形態素解析とは、文を形態素（例えば、言語で意味を持つ最小単位）の列に分割し、接続詞や助詞を取り除く。形態素解析には様々な公知の手法があるが、いずれの手法を用いてもよい。

【 0 0 4 3 】

次に、ステップ S 1 0 5 では、単語抽出部 1 3 0 2 が、ドキュメントデータのテキストを形態素解析した結果の中から、全字種の語彙を抽出する。そして、全字種の語彙を、全字種用語として、記憶部 1 5 の全字種用語リスト 1 5 0 4 に記憶する。

10

【 0 0 4 4 】

次に、ステップ S 1 0 6 では、上述のステップ S 1 0 5 を行うと共に、単語抽出部 1 3 0 2 が、ドキュメントデータのテキストを形態素解析した結果の中から、カタカナ語彙を抽出する。

【 0 0 4 5 】

次に、ステップ S 1 0 7 では、カタカナ用語抽出部 1 3 0 3 が、カタカナ語彙の用語ごとに重要度（後述）を計算し、管理者の設定する閾値以上の用語を特定する。なお、カタカナ用語特定処理の詳細については、図 4 で後述する。

【 0 0 4 6 】

次に、ステップ S 1 0 8 では、カタカナ用語抽出部 1 3 0 3 が、カタカナ語彙に対してカタカナ用語特定処理を行い特定した用語群を抽出して、記憶部 1 5 のカタカナ用語リスト 1 5 0 3 に記憶する。

20

【 0 0 4 7 】

次に、ステップ S 1 0 9 では、専門用語抽出部 1 3 0 4 が、カタカナ用語リスト 1 5 0 3 を用いて、全字種用語リスト 1 5 0 4 の用語群の中から専門用語を特定する。なお、専門用語特定処理の詳細については、図 5 で後述する。

【 0 0 4 8 】

次に、ステップ S 1 1 0 では、専門用語抽出部 1 3 0 4 が、全字種用語リスト 1 5 0 4 から専門用語特定処理を行い特定した用語群を、専門用語として抽出する。そして、抽出した専門用語と共に、カタカナ用語リスト 1 5 0 3 の用語を専門用語として、専門用語辞書に登録してもよい。

30

【 0 0 4 9 】

図 4 は、本発明の一実施形態に係るカタカナ用語特定処理のフローチャートである。

【 0 0 5 0 】

まず、ステップ S 1 7 1 では、専門用語抽出装置 1 0 の制御部 1 3 が、カタカナ用語抽出部 1 3 0 3 により、カタカナ語彙について用語ごとに重要度を計算する。なお、重要度の計算方法は、FLR (Frequency Left Right) 法、C - Value (Collocation - Value) 法、MC - Value (Modified Collocation - Value) 法などがあるので以下説明する。

40

【 0 0 5 1 】

FLR 法は、接続頻度 LR 法又は接続種類 LR 法に、用語 W がドキュメントデータ中に出現した頻度 F を加味する方法である。詳細は（非特許文献 1）を参照。接続頻度 LR 法は、語彙を走査し、用語 W を構成する単語について、該単語の左右それぞれに単語が出現する回数を計算する。又、接続種類 LR 法は、単語の左右それぞれに何種類の単語が出現するかをカウントする。ここで、例えば、カタカナ語彙中の用語「サーバシステム、コンピュータシステム、オープンシステム」があり、構成する単語を分けると（サーバ | システム）、（コンピュータ | システム）、（オープン | システム）となり、単語「システム」の左に単語が 3 回出現したので、単語「システム」の接続頻度 LR 法での左方スコアは  $L(\text{システム}) = 3$  となる。又、単語「システム」の左に単語が 3 種類出現したので、連

50

接種類 LR 法での左方スコアは  $L(\text{システム}) = 3$  となる。

【0052】

一般に、単語  $w_1, w_2, \dots, w_n$  が連なって構成する用語  $W = w_1, w_2, \dots, w_n$  について、接続頻度 LR 法又は接続種類 LR 法の用語  $W$  のスコア  $LR(W)$  が、数 1 のように定義される。

【数 1】

$$LR(W) = \left( \prod_{i=1}^n (L(w_i) + 1)(R(w_i) + 1) \right)^{\frac{1}{2n}}$$

10

$n$  : 単語数

$L(W_i), R(W_i)$  : 単語  $W_i$  の左右それぞれに単語が出現する回数又は種類数

【0053】

そして、接続頻度  $LR(W)$  又は接続種類  $LR(W)$  に、用語  $W$  がドキュメントデータ中に出現した頻度  $F(W)$  を加味した、重要度  $FLR(W)$  が、数 2 のように定義される。

【数 2】

$$FLR(W) = F(W) \times LR(W)$$

20

$F(W)$  : 用語  $W$  のドキュメントデータ中の出現頻度

$LR(W)$  : 用語  $W$  の接続頻度 LR 又は接続種類 LR

【0054】

又、C-Value 法は、用語  $W = w_1, w_2, \dots, w_n$  について、重要度 C-Value  $(W)$  が、数 3 のように定義される。C-Value 法についての詳細は (Katerina T. Frantzi and Sophia Ananiadou, Extracting nested collocations. In COLING '96, pp. 41-46, 1996.) を参照。

30

【数 3】

$$C-Value(W) = (n-1) \times \left( F(W) - \frac{T(W)}{C(W)} \right)$$

$n$  : 単語数

$T(W)$  : 用語  $W$  を部分文字列として含むより長い用語の出現頻度

$C(W)$  : 用語  $W$  を部分文字列として含むより長い用語の種類数

$F(W)$  : 用語  $W$  のドキュメントデータ中の出現頻度

40

【0055】

なお、C-Value 法は、 $n = 1$  のとき (用語が単一の単語だけからなるとき) 0 (ゼロ) になり、適切な重要度を示さない。そこで、MC-Value 法では、 $n = 1$  の場合でも重要度を計算できるよう、 $(n-1)$  の代わりに  $n$  を用いている。ここで、用語  $W = w_1, w_2, \dots, w_n$  について、重要度 MC-Value  $(W)$  が、数 4 のように定義される。MC-Value 法についての詳細は (非特許文献 1) を参照。

【数 4】

$$MC - Value(W) = n \times (F(W) - \frac{T(W)}{C(W)})$$

n : 単語数

T ( W ) : 用語 W を部分文字列として含むより長い用語の出現頻度

C ( W ) : 用語 W を部分文字列として含むより長い用語の種類数

F ( W ) : 用語 W のドキュメントデータ中の出現頻度

10

【 0 0 5 6 】

次に、ステップ S 1 7 2 では、カタカナ用語抽出部 1 3 0 3 が、カタカナ語彙から、管理者が設定した閾値以上の重要度の用語を特定する。このようにして、カタカナ語彙から、カタカナの専門用語を特定することができる。

【 0 0 5 7 】

図 5 は、本発明の一実施形態に係る専門用語特定処理のフローチャートである。

【 0 0 5 8 】

まず、ステップ S 1 9 1 では、専門用語抽出装置 1 0 の制御部 1 3 が、専門用語抽出部 1 3 0 4 により、カタカナ用語リスト 1 5 0 3 を用いて、全字種用語リスト 1 5 0 4 のそれぞれの用語について、共起ヒット情報を計算する。ここで、共起ヒット情報の計算方法は、シン普森係数値と T F ・ I D F ( T e r m F r e q u e n c y ・ I n v e r s e D o c u m e n t F r e q u e n c y ) 法とを用いる。

20

【 0 0 5 9 】

なお、シン普森係数値は、用語と用語の共起の強さを測る尺度であり、スコアが 0 ~ 1 の範囲で、高いほど共起が強い。そして、カタカナ用語 X と全字種用語 Y についての、シン普森係数値 R ( X , Y ) が、数 5 のように定義される。

【数 5】

$$R(X, Y) = \frac{|X \cap Y|}{\min(|X|, |Y|)}$$

30

| X | : カタカナ用語 X の単独ヒット数

| Y | : 全字種用語 Y の単独ヒット数

| X Y | : カタカナ用語 X と全字種用語 Y の A N D 検索でのヒット数

【 0 0 6 0 】

次に、ステップ S 1 9 2 では、専門用語抽出部 1 3 0 4 が、共起ヒット情報を基に、管理者が設定した閾値以上の共起の強さを持つ用語を、専門用語として特定する。

【 0 0 6 1 】

[ 実施例 1 ]

40

以下、専門分野として「アダルト専門分野」を例に、カタカナ用語抽出部 1 3 0 3 による、重要度 F L R の計算方法を具体的に説明する。図 6 は、実施例 1 に係るアダルト専門分野のドキュメントデータのカタカナ語彙を示す図である。図 6 ( a ) は、カタカナ語彙中の単語「セックス」を含む用語群である。図 6 ( b ) は、単語「セックス」の左右接続単語の出現頻度である。図 6 ( c ) は、カタカナ語彙中の単語「パートナー」を含む用語群である。図 6 ( d ) は、単語「パートナー」の左右接続単語の出現頻度である。ここで、単語「セックス」と単語「パートナー」について F L R 法による重要度を計算する。

【 0 0 6 2 】

まず、接続頻度法に基づき、重要度 F L R を計算する。カタカナ語彙中の単語「セックス」を含む用語群 ( 図 6 ( a ) ) において、単語「セックス」の出現頻度 F ( セックス )

50

は  $n = 3$  である。そして、図 6 ( b ) に示すように、単語「セックス」の左接単語は、「アナルセックス ( 3 )、テレホンセックス ( 1 )、オーラルセックス ( 1 )」であることから、頻度  $L$  (セックス) が  $i = 5$  となる。又、右接単語は、「セックスパートナー ( 2 )、セックスレス ( 1 )」であることから、頻度  $R$  (セックス) が  $i = 3$  となる。ここで、接続頻度に基づく重要度  $FLR$  (セックス) を計算する。

【数 6】

$$\begin{aligned}
 FLR(\text{セックス}) &= F(\text{セックス}) \times LR(\text{セックス}) \\
 &= F(\text{セックス}) \times \prod_{i=1}^1 (L(\text{セックス})+1)(R(\text{セックス})+1)^{\frac{1}{2}} \\
 &= F(\text{セックス}) \times \sqrt{(L(\text{セックス})+1)(R(\text{セックス})+1)} \\
 &= 3 \times \sqrt{(5+1)(3+1)} \\
 &= 3 \times \sqrt{24} \\
 &\approx 14.70
 \end{aligned}$$

10

このようにして、接続頻度法に基づく重要度  $FLR$  (セックス) は 14.70 と計算される。

【0063】

続いて、カタカナ語彙中の単語「パートナー」を含む用語群 ( 図 6 ( c ) ) において、単語「パートナー」の出現頻度  $F$  (パートナー) は  $n = 2$  である。そして、図 6 ( d ) に示すように、単語「パートナー」の左接単語は、「セックスパートナー ( 2 )」であることから、頻度  $L$  (セックス) が  $i = 2$  となる。又、右接単語は、「パートナーリレーション ( 1 )」であることから、頻度  $R$  (セックス) が  $i = 1$  となる。ここで、接続頻度に基づく重要度  $FLR$  (パートナー) を計算する。

20

【数 7】

$$\begin{aligned}
 FLR(\text{パートナー}) &= F(\text{パートナー}) \times LR(\text{パートナー}) \\
 &= F(\text{パートナー}) \times \prod_{i=1}^1 (L(\text{パートナー})+1)(R(\text{パートナー})+1)^{\frac{1}{2}} \\
 &= F(\text{パートナー}) \times \sqrt{(L(\text{パートナー})+1)(R(\text{パートナー})+1)} \\
 &= 2 \times \sqrt{(2+1)(1+1)} \\
 &= 2 \times \sqrt{6} \\
 &\approx 4.9
 \end{aligned}$$

30

このようにして、接続頻度法に基づく重要度  $FLR$  (パートナー) は 4.9 と計算される。

【0064】

次は、接続種類法に基づく、重要度  $FLR$  を計算する。カタカナ語彙中の単語「セックス」を含む用語群 ( 図 6 ( a ) ) において、単語「セックス」の出現頻度  $F$  (セックス) = 3 である。そして、図 6 ( b ) に示すように、単語「セックス」の左接単語は、「アナルセックス、テレホンセックス、オーラルセックス」であることから、種類  $L$  (セックス) が  $i = 3$  となる。又、右接単語は、「セックスパートナー、セックスレス」であることから、種類  $R$  (セックス) が  $i = 2$  となる。ここで、接続種類に基づく重要度  $FLR$  (セックス) を計算する。

40

## 【数 8】

$$\begin{aligned}
 FLR(\text{セックス}) &= F(\text{セックス}) \times LR(\text{セックス}) \\
 &= F(\text{セックス}) \times \prod_{i=1}^1 (L(\text{セックス})+1)(R(\text{セックス})+1)^{\frac{1}{2}} \\
 &= F(\text{セックス}) \times \sqrt{(L(\text{セックス})+1)(R(\text{セックス})+1)} \\
 &= 3 \times \sqrt{(3+1)(2+1)} \\
 &= 3 \times \sqrt{12} \\
 &\approx 10.4
 \end{aligned}$$

10

このようにして、接続種類法に基づく重要度 F L R (セックス) は 10.4 と計算される。

## 【0065】

続いて、カタカナ語彙中の単語「パートナー」を含む用語群(図6(c))において、単語「パートナー」の出現頻度 F (パートナー) = 2 である。そして、図6(d)に示すように、単語「パートナー」の左接単語は、「セックスパートナー」であることから、種類 L (セックス) が i = 1 となる。又、右接単語は、「パートナーリレーション」であることから、種類 R (セックス) が i = 1 となる。ここで、接続種類法に基づく重要度 F L R (パートナー) を計算する。

20

## 【数 9】

$$\begin{aligned}
 FLR(\text{パートナー}) &= F(\text{パートナー}) \times LR(\text{パートナー}) \\
 &= F(\text{パートナー}) \times \prod_{i=1}^1 (L(\text{パートナー})+1)(R(\text{パートナー})+1)^{\frac{1}{2}} \\
 &= F(\text{パートナー}) \times \sqrt{(L(\text{パートナー})+1)(R(\text{パートナー})+1)} \\
 &= 2 \times \sqrt{(1+1)(1+1)} \\
 &= 3 \times \sqrt{4} \\
 &= 6
 \end{aligned}$$

30

このようにして、接続種類法に基づく重要度 F L R (パートナー) は 6 と計算される。

## 【0066】

このように、F L R 法に基づき、重要度を計算することができる。そして、閾値以上の重要度の用語を、専門用語として特定する。ここで、例えば、接続頻度において、F L R (セックス) が 14.70、F L R (パートナー) が 4.9 の場合、閾値を 8 に設定することで、単語「セックス」のみが専門用語として特定できる。又、接続種類において、F L R (セックス) が 10.4、F L R (パートナー) が 6 の場合、閾値を 8 に設定することで、単語「セックス」のみが専門用語として特定できる。こうすることにより、カタカナ語彙中の用語から、閾値以上の重要度の用語を、アダルト専門分野のカタカナの専門用語として特定できる。

40

## 【0067】

次に、専門用語抽出部 1304 による、共起ヒットの計算方法を具体的に説明する。図7は、実施例1に係る共起ヒットの具体例を示す図である。

## 【0068】

まず、カタカナ用語リスト 1503 のカタカナ用語「セックス」と、全字種用語リスト 1504 の全字種用語「胸チラ」とについて、シンプソン係数値を計算する。ここで、図7に示す、ドキュメントデータにおける、カタカナ用語「セックス」の単独ヒット数(検索して抽出された数)は 7009、全字種用語「胸チラ」の単独ヒット数は 452、カタカナ用語「セックス」と全字種用語「胸チラ」とで AND 検索したヒット数は 414 であ

50

る。ここで、シン普森係数値  $R$  (セックス, 胸チラ) が、数 10 のように計算される。

【数 10】

$$R(\text{セックス}, \text{胸チラ}) = \frac{|\text{セックス} \cap \text{胸チラ}|}{\min(|\text{セックス}|, |\text{胸チラ}|)} = 414/452 = 0.915$$

$$|\text{セックス}| = 7009$$

$$|\text{胸チラ}| = 452$$

$$|\text{セックス} \cap \text{胸チラ}| = 414$$

10

このことにより、カタカナ用語「セックス」と全字種用語「胸チラ」との共起の強さが 0.915 となり、1 に近いので共起が強いことがわかる。

【0069】

次に、カタカナ用語リスト 1503 のカタカナ用語「セックス」と、全字種用語リスト 1504 の全字種用語「週末」とについて、シン普森係数値を計算する。ここで、図 7 に示す、ドキュメントデータにおける、カタカナ用語「セックス」の単独ヒット数は 7009、全字種用語「週末」の単独ヒット数は 1063、カタカナ用語「セックス」と全字種用語「週末」とで AND 検索したヒット数は 278 である。ここで、シン普森係数値  $R$  (セックス, 週末) は数 11 のように計算される。

20

【数 11】

$$R(\text{セックス}, \text{週末}) = \frac{|\text{セックス} \cap \text{週末}|}{\min(|\text{セックス}|, |\text{週末}|)} = 278/1063 = 0.262$$

$$|\text{セックス}| = 7009$$

$$|\text{週末}| = 1063$$

$$|\text{セックス} \cap \text{週末}| = 278$$

このことにより、カタカナ用語「セックス」と全字種用語「週末」との共起の強さが 0.262 となり、0 (ゼロ) に近いので共起が弱いことがわかる。

30

【0070】

このようにして、カタカナ用語リスト 1503 のカタカナ用語と、全字種用語リスト 1504 の全字種用語とについて、シン普森係数値を計算する。そして、全字種用語リスト 1504 の全字種用語を、シン普森係数値で降順にソートし、専門用語を抽出するが、いくつかの問題点がある。ここで、シン普森係数値の問題点と解決方法とについて、図 8 に基づき説明する。

【0071】

図 8 は、実施例 1 に係る全字種用語リスト 1504 の全字種用語をシン普森係数値で降順にソートした図である。はじめの行には全字種用語「風俗店」がシン普森係数値 = 1.000 であることが示されている。同様に、シン普森係数値の降順に全字種用語が並ぶ。

40

【0072】

ここで、シン普森係数値の問題点の 1 つには、低頻度な全字種用語について、共起が少なくノイズである場合が多いが、シン普森係数値が高くなりやすい問題がある。例えば、図 8 の 6 行目 全字種用語「具体案」の 0.667 は、カタカナ用語リスト 1503 のカタカナ用語「セックス」との、シン普森係数値を示す。ここでは、ドキュメントデータにおける、カタカナ用語「セックス」の単独ヒット数は 7009、全字種用語「具体案」の単独ヒット数は 3、カタカナ用語「セックス」と全字種用語「具体案」とで AND 検索したヒット数は 2 である。ここで、シン普森係数値は数 12 のように計算されて

50

いる。

【数 1 2】

$$R(\text{セックス}, \text{具体例}) = \frac{|\text{セックス} \cap \text{具体例}|}{\min(|\text{セックス}|, |\text{具体例}|)} = 2/3 = 0.667$$

$$|\text{セックス}| = 7009$$

$$|\text{具体案}| = 3$$

$$|\text{セックス} \cap \text{具体案}| = 2$$

10

このことにより、カタカナ用語「セックス」と全字種用語「具体案」との共起の強さが 0.667 となる。しかし、この場合、カタカナ用語「セックス」の単独ヒット数が 7009 であるのに対して、全字種用語「具体案」の単独ヒット数が 3 と低頻度である。よって、共起が強いとはいえない。そこで、全字種用語の単独ヒット数について閾値を設けることで解決することができる。例えば、閾値を 4 に設定することにより、全字種用語「具体案」について全字種用語から除くことができる。

【0073】

しかし、閾値を設定することにより、どのカタカナ用語とも共起するような全字種用語（いわゆる一般語）は、単独ヒット数が多く、シン普森係数値が高くなりやすい問題がある。例えば、図 8 の 10 行目 全字種用語「フリーウェア」の 0.613 は、カタカナ用語リスト 1503 のカタカナ用語「セックス」との、シン普森係数値を示す。ここでは、ドキュメントデータにおける、カタカナ用語「セックス」の単独ヒット数は 7009、全字種用語「フリーウェア」の単独ヒット数は 62、カタカナ用語「セックス」と全字種用語「フリーウェア」とで AND 検索した単独ヒット数は 38 である。ここで、シン普森係数値は数 13 のように計算されている。

20

【数 1 3】

$$R(\text{セックス}, \text{フリーウェア}) = \frac{|\text{セックス} \cap \text{フリーウェア}|}{\min(|\text{セックス}|, |\text{フリーウェア}|)}$$

$$= 38 / 62 = 0.613$$

$$|\text{セックス}| = 7009$$

$$|\text{フリーウェア}| = 62$$

$$|\text{セックス} \cap \text{フリーウェア}| = 38$$

30

このことにより、カタカナ用語「セックス」と全字種用語「フリーウェア」との共起の強さが 0.613 となる。しかし、全字種用語「フリーウェア」は一般語であるので、全字種用語から除く。そこで、閾値を 63 に設定することにより、全字種用語「フリーウェア」について全字種用語から除くことができるが、他の全字種用語も除かれてしまう。そこで、公知の TF・IDF 法を用いて解決をする。

40

【0074】

TF・IDF 法は、ドキュメントの特徴を示す単語を抽出する方法であり、ドキュメントデータの特定のページに偏って多く出現する単語ほど高スコアとなる。なお、単語 X についての TF・IDF 値は、数 14 のように定義される。

【数 1 4】

$$TF \cdot IDF = TF \cdot \log\left(\frac{N}{DF}\right)$$

50

TF：単語 X の全ページ中の出現頻度

DF：単語 X のページ頻度（いくつのページに跨って出現したか）

N：総ページ数

ここで、具体的な例を示す。

#### 【0075】

まず全字種用語「胸チラ」が、TF値 = 1423、IDF値 = 6.059である場合、TF・IDF値は8622.953となる。又、全字種用語「フリーウェア」が、TF値 = 97、IDF値 = 7.799である場合、TF・IDF値は756.542となる。ここで、全字種用語「胸チラ」はTF・IDF値が高くドキュメントデータの特定のページに偏って多く出現していることが分かる。そして、全字種用語「フリーウェア」はTF・IDF値が低いのでドキュメントデータ全体に、一般語として使われていることが分かる。このことにより、例えば全字種用語「フリーウェア」のような、どのカタカナ用語とも共起するような全字種用語について、TF・IDF値を用い、閾値を設けることで全字種用語から除くことができる。ここで、専門用語特定処理にシンプソン係数値とTF・IDF法とを用いた具体的な例について、図9に基づき説明する。

10

#### 【0076】

図9は、実施例1に係る全字種用語リスト1504から抽出された全字種用語を示す図である。ここで、全字種用語は、第1キー：シンプソン係数値、第2キー：TF・IDF値でソートしている。そして、シンプソン係数値を求める際に、全字種用語の単独ヒット数の閾値を56に設定し、低頻度な全字種用語を除いてある。又、全字種用語のTF・IDF値について、閾値を760に設定し、どのカタカナ用語とも共起するような全字種用語を除いてある。このようにして、共起ヒット情報に基づいた、全字種用語を特定することができる。

20

#### 【0077】

ここで、専門用語特定処理後の全字種用語の具体的な例について、図10に基づき説明する。

#### 【0078】

図10は、実施例1に係る専門用語特定処理後の全字種用語を示す図である。

#### 【0079】

図10に示すように、全字種用語リスト1504の全字種用語から、専門用語として、「風俗店、女王、風俗嬢、・・・」といった全字種用語が抽出されている。又、全字種用語であったが、単独ヒット数の閾値を設定することで、「具体案、介護士、やすみ、・・・」といった全字種用語を、専門用語とすることが回避されている。さらに、TF・IDF値を用いて、「フリーウェア、行楽地、株投資、・・・」といった全字種用語を、専門用語とすることが回避されている。すなわち、抽出された専門用語は、アダルト専門分野のWebサイト20に使用されている専門用語であるため、掲載禁止用語として用いることができる。そして、抽出した掲載禁止用語と共に、カタカナ用語リスト1503の用語を掲載禁止用語として、掲載禁止用語辞書に登録してもよい。

30

#### 【0080】

なお、本発明の専門用語抽出には、形態素解析を用いるが、形態素解析後の品詞の並びを参照して、連続した単語を抽出してもよい。つまり、単独では専門用語とならない単語でも、単語同士を組み合わせた場合に、掲載禁止用語となる専門用語を抽出する。例えば、単語「女子高生」と単語「画像」は、それぞれ一般的な用語であるが、2つの単語を組み合わせた用語「女子高生画像」を、全字種用語として抽出する。そして、アダルト専門分野のWebサイト20のドキュメントデータにおいて、カタカナ用語「」などとの共起の強さを計算し、専門用語として特定する。そして、全字種用語「女子高生画像」を掲載禁止用語として抽出できる。

40

#### 【0081】

#### [実施例2]

以下、専門分野として、ロボット工学関連分野を対象とした実施例を説明する。

50

## 【 0 0 8 2 】

専門用語抽出装置 1 0 を含むシステム 1 の構成及び機能ブロックは、図 2 と同様である。ここでは、URL リスト 1 5 0 1 は、ロボット工学関連分野の URL が設定された URL リスト 1 5 0 1 d を使用する。

## 【 0 0 8 3 】

又、専門用語抽出処理における実施形態は、図 3 と同様である。ここでは、クローラ部 1 3 0 1 が、ロボット工学関連分野の URL リスト 1 5 0 1 d に基づき、ロボット工学関連分野の Web サイト 2 0 のドキュメントデータを収集し、コンテンツリポジトリ 1 5 0 2 に記憶する。そして、単語抽出部 1 3 0 2 が、コンテンツリポジトリ 1 5 0 2 のドキュメントデータを形態素解析し、カタカナ語彙と全字種の語彙を抽出する。ここで、全字種の語彙は、全字種用語リスト 1 5 0 4 に記憶する。

10

## 【 0 0 8 4 】

次に、カタカナ用語抽出部 1 3 0 3 が、カタカナ語彙の用語ごとに重要度を計算し、管理者の設定する閾値以上の用語を特定するカタカナ用語特定処理については、図 4 と同様である。そして、カタカナ語彙において特定した用語をカタカナ用語リスト 1 5 0 3 に記憶する。

## 【 0 0 8 5 】

次に、専門用語抽出部 1 3 0 4 が、カタカナ用語リスト 1 5 0 3 と、全字種用語リスト 1 5 0 4 とにおいて共起の強い全字種用語を専門用語として特定する専門用語特定処理については、図 5 と同様である。ここで、シン普森係数値と TF・IDF 法とを用いた専門用語特定処理の具体的な例について、図 1 1 に基づき説明する。

20

## 【 0 0 8 6 】

図 1 1 は、実施例 2 に係る全字種用語リスト 1 5 0 4 から抽出された全字種用語を示す図である。ここで、ロボット工学関連分野の Web サイト 2 0 の Web ページデータから抽出した全字種用語は、第 1 キー：シン普森係数値、第 2 キー：TF・IDF 値でソートしている。そして、シン普森係数値を求める際に、全字種用語の単独ヒット数の閾値を 8 に設定し、低頻度な全字種用語を除いてある。又、全字種用語の TF・IDF 値について、閾値を 1 5 に設定し、どのカタカナ用語とも共起するような全字種用語を除いてある。このようにして、共起ヒット情報に基づいて、「ロボ」、「ゲーム」、「大会」、  
・  
・  
といった、ロボット工学関連分野における全字種用語を特定している。

30

## 【 0 0 8 7 】

そして、専門用語抽出部 1 3 0 4 が、共起ヒット情報を基に、管理者が設定した閾値以上の共起の強さを持つ全字種用語を、専門用語として抽出する。ここで、専門用語特定処理後の全字種用語の具体的な例について、図 1 2 に基づき説明する。

## 【 0 0 8 8 】

図 1 2 は、実施例 2 に係る専門用語特定処理後の全字種用語を示す図である。

## 【 0 0 8 9 】

図 1 2 に示すように、全字種用語リスト 1 5 0 4 の全字種用語から、専門用語として、「ロボ、ゲーム、大会、歩行、ASIMO（登録商標）、ソニー（登録商標）・・・」といった全字種用語が抽出されている。又、全字種用語であったが、単独ヒット数の閾値を設定することで、「アリーナ、ポケモン（登録商標）、ユニーク、・・・」といった全字種用語を、専門用語とすることが回避されている。さらに、TF・IDF 値を用いて、「Copyright、TOKYO、http、・・・」といった全字種用語を、専門用語とすることが回避されている。そして、抽出した専門用語と共に、カタカナ用語リスト 1 5 0 3 の用語を専門用語として、ロボット工学関連分野の専門用語辞書に登録してもよい。さらに、専門用語をロボット工学関連分野の情報を収集するキーワードとして用いるなど、様々なことに用いることができる。

40

## 【 0 0 9 0 】

## [ 共起ヒット情報の別の計算方法 ]

以上、共起ヒット情報の計算方法は、シン普森係数値と TF・IDF 法とを用いて説

50

明したが、シン普森係数値に代えて、相互情報量値、ダイス係数値、ジャガード係数値、コサイン類似度値を用いてもよい。ここで、カタカナ用語 X と全字種用語 Y の共起の強さを示す、相互情報量値は、数 15 のように定義される。

【数 15】

$$R(X, Y) = \log \frac{N|X \cap Y|}{|X||Y|}$$

10

| X | : カタカナ用語 X の単独ヒット数

| Y | : 全字種用語 Y の単独ヒット数

| X Y | : カタカナ用語 X と全字種用語 Y との AND 検索でのヒット数

N : 総ページ数

【0091】

次に、カタカナ用語 X と全字種用語 Y の共起の強さを示す、ダイス係数値は、数 16 のように定義される。

【数 16】

$$R(X, Y) = \frac{2|X \cap Y|}{|X| + |Y|}$$

20

| X | : カタカナ用語 X の単独ヒット数

| Y | : 全字種用語 Y の単独ヒット数

| X Y | : カタカナ用語 X と全字種用語 Y との AND 検索でのヒット数

【0092】

次に、カタカナ用語 X と全字種用語 Y の共起の強さを示す、ジャガード係数値は、数 17 のように定義される。

30

【数 17】

$$R(X, Y) = \frac{|X \cap Y|}{|X \cup Y|}$$

| X Y | : カタカナ用語 X と全字種用語 Y との AND 検索でのヒット数

| X Y | : カタカナ用語 X と全字種用語 Y の OR 検索でのヒット数

40

【0093】

次に、カタカナ用語 X と全字種用語 Y の共起の強さを示す、コサイン類似度値は、数 18 のように定義される。

【数 18】

$$R(X, Y) = \frac{|X \cap Y|}{\sqrt{|X||Y|}}$$

50

| X | : カタカナ用語 X の単独ヒット数

| Y | : 全字種用語 Y の単独ヒット数

| X Y | : カタカナ用語 X と全字種用語 Y との A N D 検索でのヒット数

【 0 0 9 4 】

[ 専門用語抽出装置のハードウェア構成 ]

図 1 3 は、本発明の一実施形態に係る専門用語抽出装置 1 0 ( 以下、単に専門用語抽出装置と呼ぶ ) のハードウェア構成を示す図である。

【 0 0 9 5 】

専門用語抽出装置は、制御部 1 3 0 を構成する CPU ( C e n t r a l P r o c e s s i n g U n i t ) 1 3 1 ( マルチプロセッサ構成では CPU 1 3 2 など複数の CPU が追加されてもよい )、バスライン 1 0 5、通信 I / F ( I / F : インターフェイス ) 1 2 0、メインメモリ 1 7 0、BIOS ( B a s i c I n p u t O u t p u t S y s t e m ) 1 8 0、USBポート 1 9 0、I / Oコントローラ 1 6 0、キーボード及びマウスなどの入力装置 1 1 0、並びに表示装置 1 4 0 を備える。

10

【 0 0 9 6 】

I / Oコントローラ 1 6 0 には、テープドライブ 1 5 1、ハードディスク 1 5 3、光ディスクドライブ 1 5 2、及び半導体メモリ 1 5 4 などの記憶部 1 5 0 を接続することができる。

【 0 0 9 7 】

BIOS 1 8 0 は、専門用語抽出装置の起動時に CPU 1 3 1 が実行するブートプログラムや、専門用語抽出装置のハードウェアに依存するプログラムなどを格納する。

20

【 0 0 9 8 】

ハードディスク 1 5 3 は、専門用語抽出装置として機能するための各種プログラム及び本発明の機能を実行するプログラムを記憶する。

【 0 0 9 9 】

光ディスクドライブ 1 5 2 としては、例えば、DVD - R O M ドライブ、CD - R O M ドライブ、DVD - R A M ドライブ、CD - R A M ドライブを使用することができる。この場合は各ドライブに対応した光ディスク 1 5 2 1 を使用する。光ディスク 1 5 2 1 から光ディスクドライブ 1 5 2 によりプログラム又はデータを読み取り、I / Oコントローラ 1 6 0 を介してメインメモリ 1 7 0 又はハードディスク 1 5 3 に提供することもできる。又、同様にテープドライブ 1 5 1 に対応したテープメディア 1 5 1 1 を主としてバックアップのために使用することもできる。

30

【 0 1 0 0 】

専門用語抽出装置に提供されるプログラムは、ハードディスク 1 5 3、光ディスク 1 5 2 1、又はメモリーカードなどの記録媒体に格納されて提供される。このプログラムは、I / Oコントローラ 1 6 0 を介して、記録媒体から読み出され、又は通信 I / F 1 2 0 を介してダウンロードされることによって、専門用語抽出装置にインストールされ実行されてもよい。

【 0 1 0 1 】

上述のプログラムは、内部又は外部の記憶媒体に格納されてもよい。ここで、記憶媒体としては、ハードディスク 1 5 3、光ディスク 1 5 2 1、又はメモリーカードの他に、M D などの光磁気記録媒体、テープメディア 1 5 1 1 を用いることができる。又、専用通信回線やインターネットなどの通信回線に接続されたサーバシステムに設けたハードディスク 1 5 3 又は光ディスクライブラリなどの記憶装置を記録媒体として使用し、通信ネットワーク 3 0 を介してプログラムを専門用語抽出装置に提供してもよい。

40

【 0 1 0 2 】

ここで、表示装置 1 4 0 は、ユーザによるデータの入力を受け付ける画面を表示したり、専門用語抽出装置による演算処理結果の画面を表示したりするものであり、ブラウン管表示装置 ( C R T )、液晶表示装置 ( L C D ) などのディスプレイ装置を含む。

【 0 1 0 3 】

50

ここで、入力装置 110 は、ユーザによる入力の受け付けを行うものであり、キーボード及びマウスなどにより構成してよい。

【0104】

又、通信 I/F 120 は、専門用語抽出装置を専用ネットワーク又は公共ネットワークを介して端末と接続できるようにするためのネットワーク・アダプタである。通信 I/F 120 は、モデム、ケーブル・モデム及びイーサネット（登録商標）・アダプタを含んでよい。

【0105】

以上の例は、専門用語抽出装置のハードウェア構成について主に説明したが、コンピュータに、プログラムをインストールして、そのコンピュータを専門用語抽出装置として動作させることにより上記で説明した機能を実現することもできる。従って、本発明において一実施形態として説明した専門用語抽出装置により実現される機能は、上述の方法を当該コンピュータにより実行することにより、あるいは、上述のプログラムを当該コンピュータに導入して実行することによっても実現可能である。

【0106】

以上、本発明の実施形態について説明したが、本発明は上述した実施形態に限るものではない。又、本発明の実施形態に記載された効果は、本発明から生じる最も好適な効果を列挙したに過ぎず、本発明による効果は、本発明の実施例に記載されたものに限定されるものではない。

【図面の簡単な説明】

【0107】

【図1】本発明の一実施形態に係るシステム1の全体構成を示す図である。

【図2】本発明の一実施形態に係る専門用語抽出装置10の機能ブロック図である。

【図3】本発明の一実施形態に係る専門用語抽出処理のフローチャートである。

【図4】本発明の一実施形態に係るカタカナ用語特定処理のフローチャートである。

【図5】本発明の一実施形態に係る専門用語特定処理のフローチャートである。

【図6】実施例1に係るカタカナ語彙を示す図である。

【図7】実施例1に係る共起ヒットの具体例を示す図である。

【図8】実施例1に係る全字種用語リスト1504の全字種用語をシンブソン係数値で降順にソートした図である。

【図9】実施例1に係る全字種用語リスト1504から抽出された全字種用語を示す図である。

【図10】実施例1に係る専門用語特定処理後の全字種用語を示す図である。

【図11】実施例2に係る全字種用語リスト1504から抽出された全字種用語を示す図である。

【図12】実施例2に係る専門用語特定処理後の全字種用語を示す図である。

【図13】本発明の一実施形態に係る専門用語抽出装置10のハードウェア構成を示す図である。

【符号の説明】

【0108】

- 1 システム
- 10 専門用語抽出装置
- 20 Webサイト
- 30 通信ネットワーク
- 1501 URLリスト
- 1502 コンテンツリポジトリ
- 1503 カタカナ用語リスト
- 1504 全字種用語リスト

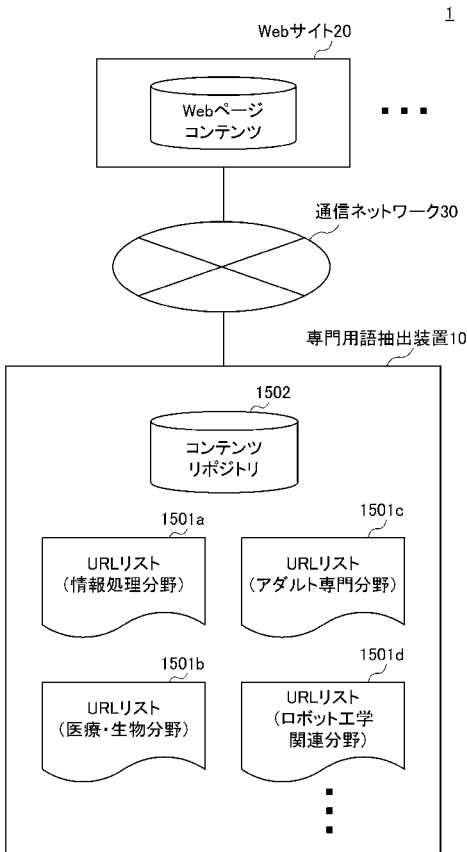
10

20

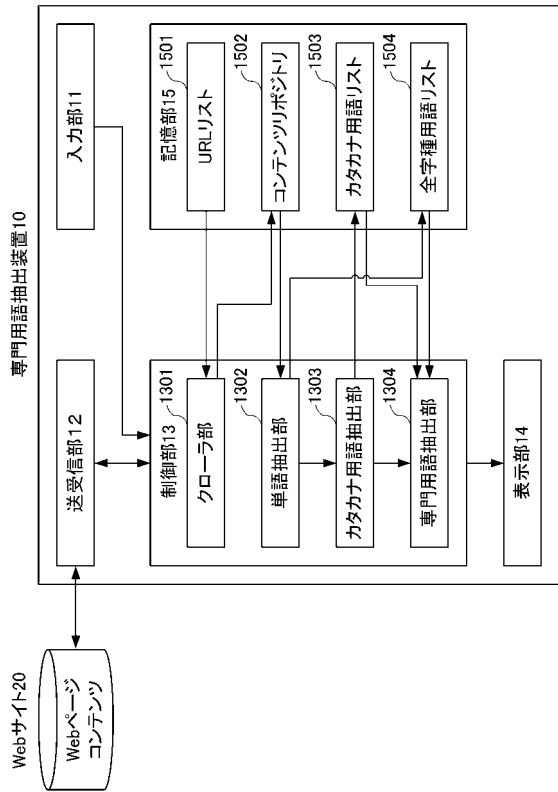
30

40

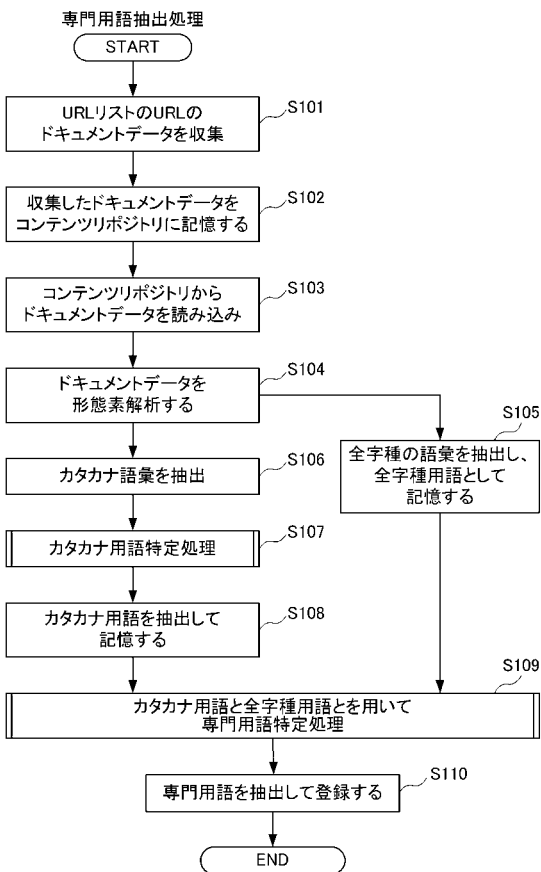
【 図 1 】



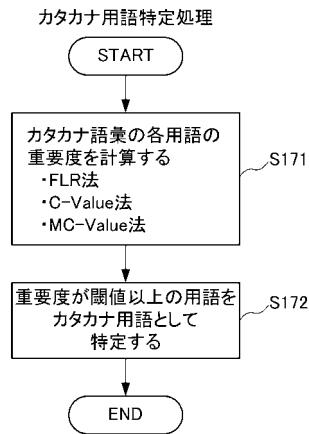
【 図 2 】



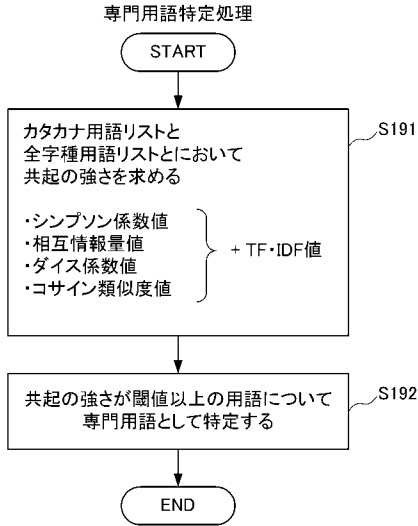
【 図 3 】



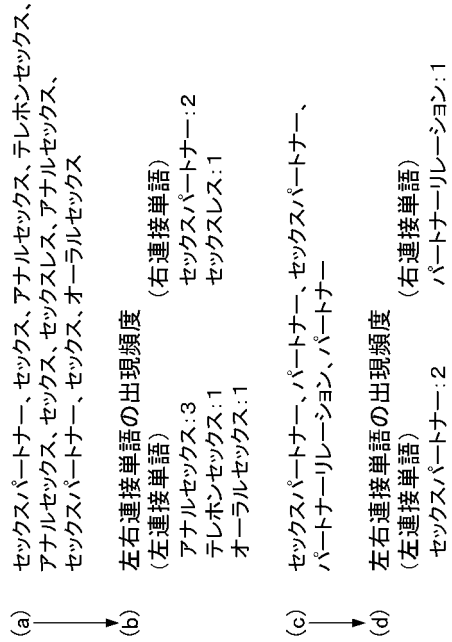
【 図 4 】



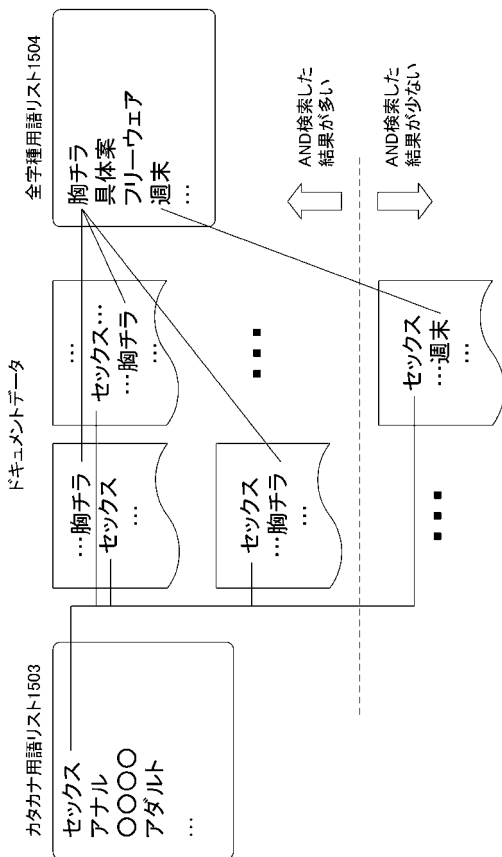
【 図 5 】



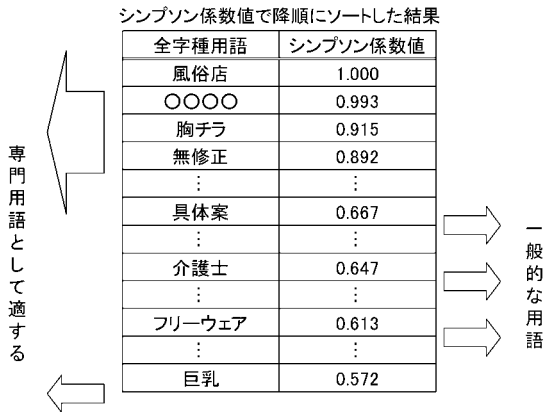
【 図 6 】



【 図 7 】



【 図 8 】



## 【 図 9 】

URLリスト:アダルト専門分野のWebサイト

全字種用語リストから抽出された全字種用語

全字種用語	シンブソン係数値	TF・IDF値
風俗店	1.000	13941.885
女王	1.000	10984.916
風俗嬢	1.000	9595.345
〇〇〇〇〇	1.000	6233.887
女性器	1.000	5084.160
同性	1.000	4477.127
処女膜	1.000	3296.157
男性器	1.000	3198.379
下着姿	1.000	2821.487
性風俗	1.000	2703.810

降順ソート

第1キー:シンブソン係数値

第2キー:TF・IDF値

閾値

単独ヒット数:56

TF・IDF値:760

## 【 図 1 0 】

アダルト専門分野のWebサイトからの専門用語(掲載禁止用語)抽出

抽出できた例
風俗店、女王、風俗嬢、〇〇〇〇〇、女性器、同性、処女膜、男性器、下着姿、性風俗、〇〇〇〇、胸チラ、無修正、巨乳
取り除けた例
●単独ヒット数の閾値で取り除けた例 具体案、介護士、やすみ、ノコギリ、東欧、搭乗客、豚トロ、鈍化、内科医、売却益、麦秋
●TF・IDF値で取り除けた例 フリーウェア、行楽地、株投資、しずく、カワイイ、撫子、茶の間、薄型、中級編、支払、超特価、浄水

## 【 図 1 1 】

URLリスト:ロボット工学関連分野のWebサイト(クロールページ数:47)

全字種用語リストから抽出された全字種用語

全字種用語	シンブソン係数値	TF・IDF値
ロボ	1.000	136.264
ゲーム	1.000	129.693
大会	1.000	101.530
歩行	1.000	100.058
ASIMO	1.000	94.215

降順ソート

第1キー:シンブソン係数値

第2キー:TF・IDF値

閾値

単独ヒット数:8

TF・IDF値:15

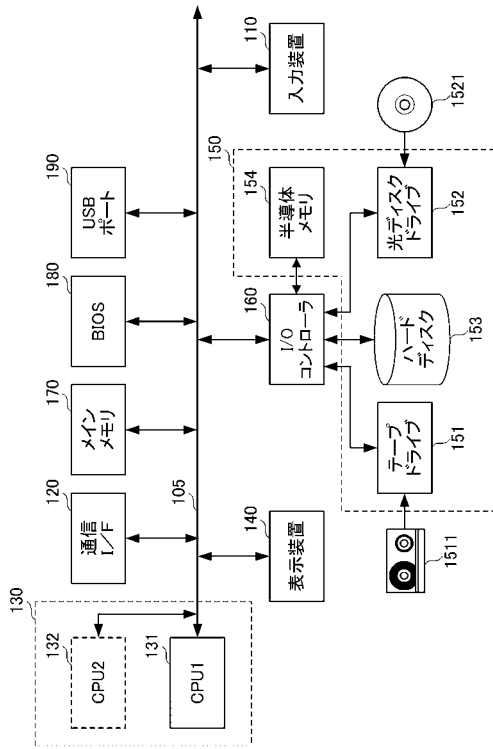
## 【 図 1 2 】

ロボット工学関連分野のWebサイトからの専門用語抽出

抽出できた例
ロボ、ゲーム、大会、歩行、ASIMO、ソニー、設定、デモ、デザイン、センサー、会話、PC、人型
取り除けた例
●単独ヒット数の閾値で取り除けた例 アリーナ、ポケモン、ユニーク、意味、完全、玩具、場合、新世界、同様、様子、連載
●TF・IDF値で取り除けた例 Copyright、TOKYO、http、ご紹介、オススメ、フロント、異なる、凶悪犯、超能力、不正、平和

【図 13】

10



フロントページの続き

Fターム(参考) 5B091 AA15 AB04 AB08 CA02 CC05 CC16